



Universiteit  
Leiden  
Governance and Global Affairs



# AI-Driven Text Analysis in the Political Economy of Sustainability 🌿

## Hybrid Retrieval-Augmented Generation and LLM Multi-Agent Approach

---

**Bastián González-Bustamante**

Leiden University

✉ [b.a.gonzalezbustamante@fgga.leidenuniv.nl](mailto:b.a.gonzalezbustamante@fgga.leidenuniv.nl)

**Natascha van der Zwan**

Leiden University

✉ [n.a.j.van.der.zwan@fgga.leidenuniv.nl](mailto:n.a.j.van.der.zwan@fgga.leidenuniv.nl)

Presentation at the ECPR Joint Sessions of Workshops, Charles University  
Prague, Czech Republic, May 20-23, 2025

# Introduction

---

# Challenges in Situating Sustainable Finance

## ➤ **Conceptual Challenge**

How to place sustainable finance within existing politico-economic typologies

## ➤ **Empirical Challenge**

Lack of official SI data and consensus on sustainable investment definition

## ➤ **Greenwashing Concern**

Difficulty in distinguishing genuine sustainable practices from greenwashing

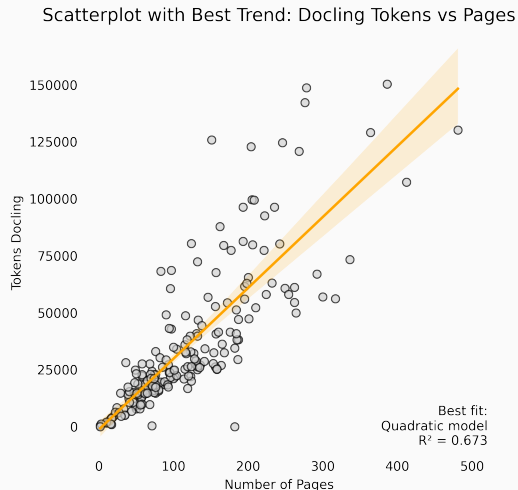
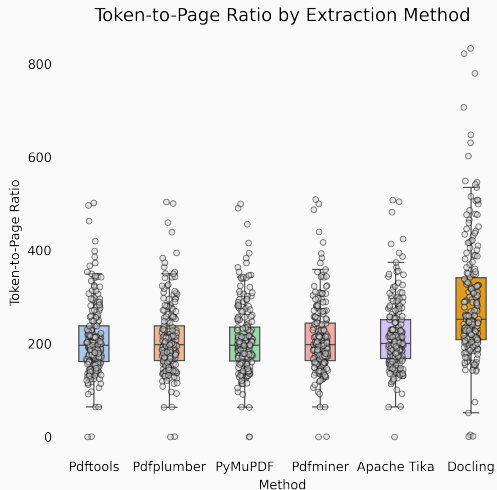


Artwork by DALL·E model

# Research Design

---

# Benchmark of Annual Reports Parsing Methods



*Note.* Two observations were excluded from the plots because of their high number of tokens.

**32**

**Countries**

Top asset owners from  
these countries  
for the moment

**24K**

**Pages**

Only annual reports,  
but we are including  
specific reports  
and disclosures

**7.6M**

**Tokens**

Tokens for process-  
ing and analysis  
for the moment

## 96 LLMs

CAP's major policy topics environment  
and energy in English for

- (1) **error rate analysis**
- (2) **orchestration model selection**

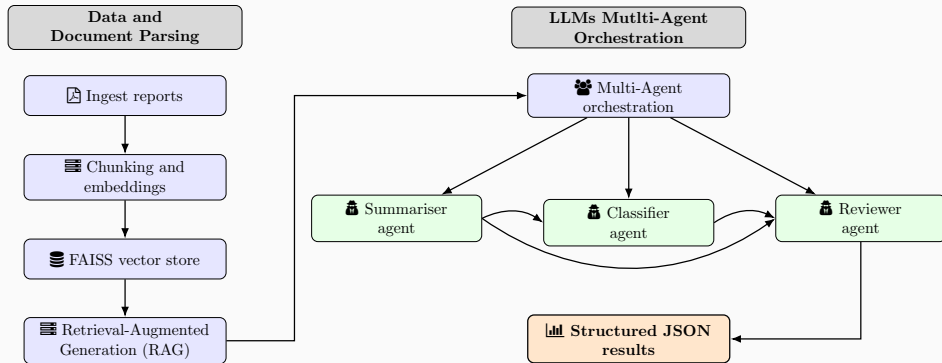
### **SOTA closed-source LLMs**

o3, o4-mini, GPT-4.1, GPT-4.5-preview,  
Gemini 2.5, Grok 3 Beta, Claude 3.7 Sonnet

### **SOTA open-source LLMs**

Llama 4 Maverick (400B) and Scout (107B),  
Mistral 3.1 (24B), Llama 3.3 (70B),  
DeepSeek-R1 (671B), DeepSeek-V3 (671B)

# RAG and Multi-Agent Orchestration Pipeline







## Summariser Prompt

You are an expert financial analyst specialising in sustainable finance disclosures. Using only the retrieved context, write a concise (200-300 word) neutral summary. Focus on these six analytical frames: emissions reporting, climate leadership, sustainable investments, green finance, climate-related risks, energy transition.

After the prose summary. Estimate how prevalent each frame is as a percentage of total sentences in the underlying document (0-100, rounded). Output a single JSON object named `frame_prevalence`, with exactly these six keys, each mapped to its % value (integer). Nothing else may follow that JSON.

## Classifier Prompt

You are an ESG taxonomy specialist. Determine how prominently the document frames every topic below, using both the prose summary and the frame\_prevalence JSON (sentence-share %): emissions reporting, climate leadership, sustainable investments, green finance, climate-related risks, energy transition.

Label each frame with one of: Not Mentioned ( $< 1\%$ ), Peripheral ( $\sim 1\text{-}5\%$ ), Moderate ( $\sim 6\text{-}15\%$ ), Central ( $> 15\%$ ).

Respond only with JSON having exactly these six keys: emissions\_reporting, climate\_leadership, sustainable\_investments, green\_finance, climate\_related\_risks, energy\_transition.



## Reviewer Prompt

You are the senior reviewer. Evaluate the summary, the frame\_prevalence JSON, and the classification below. If anything is missing, mislabelled, inconsistent, or exceeds the 300-word limit, produce corrected versions.

Return only a JSON object with three keys: revised\_summary (a refined summary ( $\leq 300$  words)), revised\_frame\_prevalence (the prevalence JSON in the same format), revised\_classification (the six-key classification JSON).

## Preliminary Results

---

Accuracy

Precision

Recall

F1-Score

o3 (2025-04-16)  
 Hermes 3 (100)  
 GPT-4o (2024-05-13)  
 Grok 3 Mini Beta  
 GPT-4.5 preview (2025-02-21)  
 Nemotron (100)  
 Grok 3 Mini Next Beta  
 Grok Beta  
 GPT-4.1 (2025-04-14)  
 GPT-4.1 mini (2025-04-14)  
 Grok 2 (12/22)  
 DeepSeek-V3 (10/18)  
 Qwen 2.5 (12/28)  
 Llama 3.1 (04/20)  
 Grok 3 Beta  
 Llama 3.1 (100)  
 GPT-4o (2024-08-06)  
 DeepSeek-V3 (01/15)  
 o3-mini (2025-01-31)  
 o1-mini (2025-04-16)  
 Grok 3 Next Beta  
 Llama 3.1 (100)  
 Gemini 2.5 Pro (03/25)  
 GPT-4o (2024-11-20)  
 Qwen 2.5 (12/28)  
 Expanse 3.5 (12/28)  
 Gemini 1.5 Pro  
 Mistral Large (24/11)  
 Athena V2 (12/28)  
 Nova (100)  
 o1-mini (2024-09-12)  
 g1 (2024-12-11)  
 Gemini 2.0 Flash  
 Gemini 2.0 Flash Exp  
 Qwen 3 (22/05)  
 Phi-4 (14/05)  
 OpenTranzer (12/28)  
 Llama 3.5 Moondance (02/05)  
 Gemini 1.5 (12/28)  
 Gemini 2.0 Flash-Lite (01/11)  
 Gemini 2.0 Flash-Lite (02-05)  
 Gemini 1.5 Flash  
 Yi Large  
 Gemini 2 Pro  
 Nova Hermes 2 (11/18)  
 Gemini 1.5 (12/28)  
 Open Mistral 8x22B  
 GPT-4o mini (2024-07-18)  
 Llama 4 Scout (10/19)  
 Mistral Small 1.1  
 Patral L2B (24/09)  
 Mistral Large (24/11)  
 Expanse 3.5 (08/28)  
 Gemini 2 (12/28)  
 Gemini 1.5 Flash (08/28)  
 Qwen 2.5 (14/05)  
 Aya Expanse (12/28)  
 G1 M4 (06/28)  
 Sonnet 2 (08/28)  
 Command R3B Max (12/28)  
 Llama 3.1 (08/28)  
 Phi 1.5 (09/28)  
 Aya (15/28)  
 Mistral Small (12/28)  
 Claude 3.5 Sonnet (2024-10-22)  
 Claude 3.2 Sonnet (2024-10-19)  
 Claude 3.5 Haiku (2024-10-22)  
 Mistral Large  
 DeepSeek-R1-D1-Llama (08/13)  
 R 1.5 (04/08)  
 GPT-4.1 nano (2025-04-14)  
 Nemotron (100)  
 Mistral OpenOrca (12/28)  
 Nemotron (100)  
 Mistral (12/28)  
 Mistral (12/28)  
 OpenTranzer (12/28)  
 Qwen 2.5 (12/28)  
 Phi-4-mini (13/08)  
 Llama 3.2 (13/08)  
 Granite 3.1 (08/28)  
 Mistral 8x22B (24/11)  
 Mistral Nano (12/28)  
 Q1 M4 (12/28)  
 Dolphin 3.0 (08/28)  
 Aya Expanse (08/28)  
 V 1.5 (08/28)  
 Gemini 3.2 (08/28)  
 Nova Hermes 2 Mistral (12/28)  
 Gemini (14/08)  
 Nemotron Mini (04/08)  
 D1 M4 (12/28)  
 D1 M4 (12/28)  
 D1 M4 (12/28)  
 Granite 3.1 Mistral (08/28)  
 Phi-1 Medium (14/08)

0.0 0.5 1.0 0.0

0.5 1.0 0.0

0.5 1.0 0.0

0.5 1.0

## Per-Frame Metrics (35% of the corpus)

Frame	Agent	Accuracy	Weighted Precision	Weighted Recall	Weighted F1-score
Emissions reporting	Classifier	0.986	0.986	0.986	0.985
Climate leadership	Classifier	0.929	0.980	0.929	0.943
Emissions reporting	Summariser	0.929	0.961	0.929	0.939
Climate leadership	Summariser	0.900	0.943	0.900	0.915
Energy transition	Classifier	0.843	0.841	0.843	0.822
Green finance	Classifier	0.800	0.809	0.800	0.797
Sustainable investments	Summariser	0.800	0.886	0.800	0.793
Energy transition	Summariser	0.786	0.809	0.786	0.792
Green finance	Summariser	0.771	0.793	0.771	0.766
Climate-related risks	Summariser	0.757	0.819	0.757	0.744
Climate-related risks	Classifier	0.686	0.749	0.686	0.655
Sustainable investments	Classifier	0.629	0.820	0.629	0.648

*Note.* Reviewer agent as ground-truth.

## Takeaways and Next Steps

---

# Takeaways and Next Steps

- ✓ **Goof Performance.** The summariser reaches 0.824 accuracy and 0.826 weighted F1-score against the reviewer yardstick. The classifier follows a fraction behind, leaving a residual gap of 1.2 points.
- ✓ **Hybrid Design.** RAG backbone, orchestration and mixed licensing strategy. Swapping in components is possible.
- ✓ **Next Steps.** There are a number of immediate improvements and next steps:
  - **Some Problems.** (1) Most of that shortfall is concentrated in frames as sustainable investments and green finance. (2) Closed-source summariser and reproducibility.
  - **Improvements.** (1) Classifier prompt needs to tie salience labels more tightly. (2) Temperature and nucleus-samples sweeps can probe whether a touch of controlled randomness. (3) Human-in-the-loop.



Thank you very much!

Do you have any questions?

[b.a.gonzalez.bustamante@fgga.leidenuniv.nl](mailto:b.a.gonzalez.bustamante@fgga.leidenuniv.nl)

[bgonzalezbustamante.com](http://bgonzalezbustamante.com)

Presentation compiled with  $\text{\LaTeX}$  and some 

 Download the latest version from [GitHub](https://github.com)